
datanator documentation

Release 0.2

Karr Lab

Oct 08, 2020

Contents

1	Contents	3
1.1	Overview	3
1.1.1	Motivation	3
1.1.2	Methodology	4
1.1.3	Supported data types and data sources	4
1.1.4	Advantages of <code>datanator</code> versus manual data aggregation	5
1.2	Installation	5
1.2.1	Install dependencies	5
1.2.2	Install <code>datanator</code>	6
1.2.3	Run <code>datanator</code>	6
1.3	Tutorial - Getting Familiar With <code>datanator</code> 's Command Line Interface	6
1.3.1	Generate Template Doc	6
1.3.2	Collecting Kinetic Data	7
1.3.3	Entering Reactions Into <code>datanator</code>	7
1.3.4	Set Maximum Proximity Limit	9
1.4	About	9
1.4.1	License	9
1.4.2	Development team	10
1.4.3	Acknowledgements	10
1.4.4	Questions and comments	10
1.5	References	10

`datanator` is a software tool for finding experimental data for building and calibrating dynamical models of cellular biochemistry such as metabolite, RNA, and protein abundances; protein complex compositions; transcription factor binding motifs; and kinetic parameters. `datanator` is particularly useful for building large models, such as whole-cell models, that require large amounts of data to constrain large numbers of parameters. `datanator` was motivated by the need for large amounts of data to constrain whole-cell models and the fact that this data is hard to utilize because it is scattered across numerous siloed repositories.

`datanator` currently supports the following data types and data sources:

- Metabolite concentrations: [ECMDB](#) and [YMBD](#)
- RNA abundance: [ArrayExpress](#)
- Protein abundance: [PaxDb](#)
- Protein complex composition: [CORUM](#)
- Transcription factor binding motifs: [JASPAR](#)
- Reaction kinetics: [SABIO-RK](#)
- Taxonomy: [NCBI Taxonomy](#)

`datanator` (1) downloads these repositories; (2) normalizes their data to a common ontology and units; (3) stores their data to a local SQLite database; and (4) provides a Python API for (a) finding relevant data to model a specific organism and environmental condition from similar species, reactions, genotypes (taxon, variant), and environments (temperature, pH, media), and (b) reducing multiple relevant observations to a single consensus recommended parameter value, and (c) exporting these consensus recommendations and their provenance to an Excel workbook. To make `datanator` easier to use, we plan to develop user-friendly command line and web-based interfaces for finding data for SBML-encoded models.

`datanator` is under active development and is not yet ready for end users. Please check back soon for updates.

This website contains detailed documentation of the `datanator` source code. Going forward, this website will also contain detailed instructions and tutorials on how to use `datanator`.

1.1 Overview

`datanator` is a software tool for finding experimental data for building and calibrating dynamical models of cellular biochemistry such as metabolite, RNA, and protein abundances; protein complex compositions; transcription factor binding motifs; and kinetic parameters. `datanator` is particularly useful for building large models such as whole-cell models that require large amounts of data to constrain large numbers of parameters.

1.1.1 Motivation

Large models such as whole-cell models are needed to help researchers, bioengineers, and physicians predict how genotype and the environment determine phenotype. In particular, large models that represent the function of each individual gene are needed to help bioengineers rationally design microorganisms to perform specific functions such as producing drugs and neutralizing pathogens and to help physicians interpret personal genomes and design personalized therapies tailored to each patient's unique genome.

Despite their potential, large models are hard to build, in part, because it is hard to aggregate the large amount of data needed to build large models because this data is scattered across numerous siloed repositories. Without data aggregation tools such as `datanator`, this data must be manually aggregated:

1. Modelers must identify appropriate data sources for their model.
2. Modelers must identify the relevant subset of this data for their model.
3. Modelers must merge inconsistently annotated data from multiple sources.
4. Modelers must reduce multiple experimental observations to individual consensus parameter value recommendations.
5. Modelers must record the provenance of each recommendation so their model is comprehensible and reproducible.

This requires extensive time and effort, is unscalable, introduces substantial selection bias into the data underlying models, and is irreproducible.

To address these problems, we have developed `datanator` to systemize and accelerate the aggregation of data for biomodels.

1.1.2 Methodology

1. `datanator` downloads data from the data repositories listed below.
2. `datanator` parses this data, normalizes this data to a common schema and common identifiers, and converts this data to common units.
3. `datanator` stores the normalized data and its provenance to a local SQLite database.
4. `datanator` helps researchers find relevant data to model a specific organism and environmental condition from similar species, reactions, genotypes, and environments according to the following filters:
 - Chemical similarity: `datanator` helps researchers identify data observed for (a) molecularly-similar species as determined by the Tanimoto similarity of their molecular fingerprints and (b) chemically-similar reactions that involve similar reaction centers and reaction mechanisms.
 - Taxonomic similarity: `datanator` helps researchers identify data observed for taxonomically-close taxa and similar genetic variants.
 - Environmental similarity: `datanator` helps researchers identify data observed for similar temperatures, pHs, and growth media.
5. `datanator` reduces multiple relevant observations to a single consensus recommended parameter value.
6. `datanator` records the provenance of the data underlying each consensus recommendation.
7. `datanator` exports the consensus recommendations and their provenance to Excel workbooks.

Currently, `datanator` provides researchers a Python API to programmatically aggregate data for models. To make `datanator` easier to use, we plan to develop user-friendly command line and web-based interfaces which will help users find data for SBML-encoded models, review this data, and generate consensus recommendations for parameter values.

1.1.3 Supported data types and data sources

`datanator` currently supports the following data type and data sources. We are also actively integrating additional data types and data sources into `datanator`.

- Databases
 - Metabolite concentrations: [ECMDB](#) and [YMBD](#)
 - RNA abundance: [ArrayExpress](#)
 - Protein abundance: [PaxDb](#)
 - Protein complex composition: [CORUM](#)
 - Transcription factor binding motifs: [JASPAR](#)
 - Reaction kinetics: [SABIO-RK](#)
 - Taxonomy: [NCBI Taxonomy](#)
- Prediction tools
 - Metabolite properties (charge, pK_a, protonation): [Open Babel](#)
 - EC number: [E-zyme](#)

1.1.4 Advantages of datanator versus manual data aggregation

- `datanator` accelerates data aggregation by merging data from multiple repositories into a single location.
- `datanator` reduces selection bias in the data used to build models by systemizing how researchers find data for models.
- `datanator` helps researchers find more data for models by helping researchers find data from chemically-similar species and reactions, taxonomically-similar organisms, and similar environmental conditions based on their temperature, pH, and media.
- `datanator` increases the comprehensibility and reproducibility of models by automatically tracking the provenance of each recommended parameter value.

1.2 Installation

The following instructions describe how to install `datanator` onto a Debian-based Linux OS using the docker image `wc_env` <https://hub.docker.com/r/karrlab/wc_env>

1.2.1 Install dependencies

First, please install the following dependencies:

- Docker
- Docker-compose

The following shell commands can be used to install these dependencies onto Ubuntu Linux 16.04:

```
apt-get update

apt-get install \
  apt-transport-https \
  ca-certificates \
  curl \
  gnupg-agent \
  software-properties-common

curl -fsSL https://download.docker.com/linux/ubuntu/gpg | sudo apt-key add -

sudo add-apt-repository \
  "deb [arch=amd64] https://download.docker.com/linux/ubuntu \
  $(lsb_release -cs) \
  stable"

sudo apt-get update

sudo apt-get install docker-ce docker-ce-cli containerd.io

sudo curl -L "https://github.com/docker/compose/releases/download/1.26.2/docker-
compose-$(uname -s)-$(uname -m)" -o /usr/local/bin/docker-compose

sudo chmod +x /usr/local/bin/docker-compose

sudo ln -s /usr/local/bin/docker-compose /usr/bin/docker-compose
```

1.2.2 Install datanator

Second, please run the following shell commands to clone and install datanator from GitHub:

```
mkdir karr_lab
mkdir ~/.wc
cd ./karr_lab
git clone git@github.com:KarrLab/pkg_utils.git
git clone git@github.com:KarrLab/wc_utils.git
git clone git@github.com:KarrLab/karr_lab_aws_manager.git
git clone git@github.com:KarrLab/datanator_query_python.git
git clone git@github.com:KarrLab/datanator.git
cd ./datanator
nano docker-compose.yml # change ``z1`` on line 13 to the proper username. Save and
↵exit by pressing ``Ctrl + X`` followed by ``Y``
docker-compose up -d
```

1.2.3 Run datanator

One needs to find the docker container ID in order to use Datanator package:

```
docker ps
docker exec -it <container_id> bash
cd karr_lab/datanator
```

All python scripts in datanator directory can be run with python3, for example:

```
python3 datanator/data_source/corum_nosql.py
```

Running the command above will parse *Corum* <<http://mips.helmholtz-muenchen.de/corum/>> and store the parsed data in KarrLab's MongoDB.

Contact [Yang](#) for any questions regarding installation and running the package.

1.3 Tutorial - Getting Familiar With datanator's Command Line Interface

1.3.1 Generate Template Doc

datanator can be used run directly from the command line.

Let's do an example

Create a new directory

In that directory run the following command:

```
$ datanator generate-template
```

This will output a template document to your directory. This template is an example of an excel sheet that can be entered into datanator

1.3.2 Collecting Kinetic Data

In order to collect kinetic information on the reaction in the template, we will use “find-kinetics”

There are three required arguments necessary to run the data collection method “find-kinetics”

1. A path to an excel sheet with the reactions you would like to search (we will discuss the necessary formatting later)
2. A name for the results file (the results are saved as an excel file, so the name should end in ‘.xlsx’)
3. The name of the species you are searching for

The first argument is a path to an excel sheet, we will use the template document we just generated labelled InputTemplate.xlsx

The second argument is a name for the output file. We will use *Results.xlsx*, but you can name it anything.

The third argument is the name of the species we are searching for. We will use ‘homo sapiens’

Run “get-kinetics”:

```
$ datanator get-kinetics InputTemplate.xlsx Results.xlsx 'homo sapiens'
```

```
$ datanator get-taxonomic-lineage 'Escherichia coli'
```

1.3.3 Entering Reactions Into datanator

Open up the template document InputTemplate.xlsx

This document has two worksheets: Reactions, and Metabolites.

Reactions Worksheet

Reaction ID	Stoichiometry
A Reaction Name #1 (ATP + UMP <=> UDP + ADP)	ATP + UMP <=> UDP + ADP
A Reaction Name #2 (GMP + ATP <=> ADP + GDP)	GMP + ATP ==> ADP + GDP

Look at the Reactions worksheet. There are two columns. The first is an identifier for the reaction. You can name this whatever you like as long as it is unique. A good suggestion is just to use the reaction string (in the template document, we did not follow this practice in order to illustrate that the name can be anything). The second column is a the stoichiometric string.

Each metabolite used in the stoichiometric string needs to be structurally defined in the “Metabolites” worksheet. Open up the “Metabolites” worksheet

Compound ID	Structure
ATP	<chem>NC1=C2N=CN(C3OC(COP([O-])(=O)OP([O-])(=O)OP([O-])([O-])=O)C(O)C3O)C2=NC=N1</chem>
UMP	<chem>OC1C(O)C(OC1COP([O-])([O-])=O)N1C=CC(=O)NC1=O</chem>
UDP	<chem>OC1C(O)C(OC1COP([O-])(=O)OP([O-])([O-])=O)N1C=CC(=O)NC1=O</chem>
AMP	<chem>NC1=C2N=CN(C3OC(COP([O-])([O-])=O)C(O)C3O)C2=NC=N1</chem>
GMP	<chem>NC1=NC2=C(N=CN2C2OC(COP([O-])([O-])=O)C(O)C2O)C(=O)N1</chem>
GDP	<chem>NC1=NC2=C(N=CN2C2OC(COP([O-])(=O)OP([O-])([O-])=O)C(O)C2O)C(=O)N1</chem>

The compound ID corresponds to the ID of the compound in the stoichiometric string. The structure is either an Inchi or a SMILES string (SMILES in this example).

Let's try adding a reaction:

Lets say we wanted to add the reaction:

```
Adenosine 3',5'-bisphosphate + H2O ==> phosphate + AMP
```

The first step is to give each metabolite in the stoichiometric string and ID without spaces. So let's change it to:

```
A-3-5-bisphosphate + H2O ==> phosphate + AMP
```

The second step is to add this reaction string to the second column in "Reactions" worksheet. In the first column, give the reaction some distinct name.

Reaction ID	Stoichiometry
A Reaction Name #1 (ATP + UMP <=> UDP + ADP)	ATP + UMP <=> UDP + ADP
A Reaction Name #2 (GMP + ATP <=> ADP + GDP)	GMP + ATP ==> ADP + GDP
A distinct name of your choosing	A-3-5-bisphosphate + H2O ==> phosphate + AMP

The third step is to structurally define each metabolite in the reaction string. We already defined structurally defined AMP (we used it in the previous reactions), so we will have to structurally define A-3-5-bisphosphate, H2O, and phosphate.

The structural information is here:

```
A-3-5-bisphosphate - NC1=C2N=CN(C3OC(COP([O-])([O-])=O)C(OP([O-])([O-])=O)C3O)C2=NC=N1
H2O - O
phosphate - OP([O-])([O-])=O
```

Now we need to add this information to the "Metabolites" worksheet.

Open up the "Metabolites" worksheet. Add the name of the compound used in the stoichiometric string (ex: phosphate) to the first column, and add the structure in the second.

Compound ID	Structure
ATP	NC1=C2N=CN(C3OC(COP([O-])([O-])=O)OP([O-])([O-])=O)C(OP([O-])([O-])=O)C3O)C2=NC=N1
UMP	OC1C(O)C(OC1COP([O-])([O-])=O)N1C=CC(=O)NC1=O
UDP	OC1C(O)C(OC1COP([O-])([O-])=O)OP([O-])([O-])=O)N1C=CC(=O)NC1=O
AMP	NC1=C2N=CN(C3OC(COP([O-])([O-])=O)C(O)C3O)C2=NC=N1
GMP	NC1=NC2=C(N=CN2C2OC(COP([O-])([O-])=O)C(O)C2O)C(=O)N1
GDP	NC1=NC2=C(N=CN2C2OC(COP([O-])([O-])=O)OP([O-])([O-])=O)C(O)C2O)C(=O)N1
A-3-5-bisphosphate	NC1=C2N=CN(C3OC(COP([O-])([O-])=O)C(OP([O-])([O-])=O)C3O)C2=NC=N1
H2O	O
phosphate	OP([O-])([O-])=O

Once again, run "get-kinetics":

```
$ datanator get-kinetics InputTemplate.xlsx Results.xlsx 'Escherichia coli'
```

You should see results for the new reaction you inputted.

1.3.4 Set Maximum Proximity Limit

The ideal kinetic data is information about the reaction being studied, collected from an experiment done in the species you are studying. However, often you will have to rely on kinetic data from different species. At a certain taxonomic distance, you might decide it's better to collect data from similar reactions taken from experiments in more closely related organisms.

There are two dimensions of granularity - reaction variation and species variation - and the user can decide which data is preferred.

The user can do this by setting the "proximLimit"

Let's try an example:

First, you want to get taxonomic information about the organism you are studying. Run:

```
$ datanator get-taxonomic-lineage 'Escherichia coli'
```

You should see:

```
1: Escherichia coli
2: Escherichia
3: Enterobacteriaceae
4: Enterobacterales
5: Gammaproteobacteria
6: Proteobacteria
7: Bacteria
8: cellular organisms
9: root
```

This is the number of nodes as you start from your organism, and climb up to the top of the taxonomic tree. Each number corresponds to a node.

So let's say you are studying *Escherichia coli*. Maybe you think that anything outside the phylum proteobacteria is too distantly related to be useful. In that case, you will want to run the "get-kinetics" argument, with an optional argument --proxim-limit. The number given after --proxim-limit is the highest node that will be considered useful. Since we have chosen Proteobacteria, that number is 6

So, run:

```
datanator get-kinetics InputTemplate.xlsx Results.xlsx 'Escherichia coli' --proxim-
-limit 6
```

1.4 About

1.4.1 License

The software is released under the MIT license:

The MIT License (MIT)

Copyright (c) 2017-2018 Karr Lab

Permission is hereby granted, free of charge, to any person obtaining a copy of this software and associated documentation files (the "Software"), to deal in the Software without restriction, including without limitation the rights

(continues on next page)

(continued from previous page)

to use, copy, modify, merge, publish, distribute, sublicense, and/or sell copies of the Software, and to permit persons to whom the Software is furnished to do so, subject to the following conditions:

The above copyright notice and this permission notice shall be included in all copies or substantial portions of the Software.

THE SOFTWARE IS PROVIDED "AS IS", WITHOUT WARRANTY OF ANY KIND, EXPRESS OR IMPLIED, INCLUDING BUT NOT LIMITED TO THE WARRANTIES OF MERCHANTABILITY, FITNESS FOR A PARTICULAR PURPOSE AND NONINFRINGEMENT. IN NO EVENT SHALL THE AUTHORS OR COPYRIGHT HOLDERS BE LIABLE FOR ANY CLAIM, DAMAGES OR OTHER LIABILITY, WHETHER IN AN ACTION OF CONTRACT, TORT OR OTHERWISE, ARISING FROM, OUT OF OR IN CONNECTION WITH THE SOFTWARE OR THE USE OR OTHER DEALINGS IN THE SOFTWARE.

1.4.2 Development team

This package was developed by the following researchers in the [Karr Lab](#) at the Icahn School of Medicine at Mount Sinai in New York, USA:

- [Yosef Roth](#)
- [Saahith Pochiraju](#)
- [Balazs Szigeti](#)
- [Jonathan Karr](#)

1.4.3 Acknowledgements

This work was supported by a National Institute of Health MIRA award [grant number 1 R35 GM 119771-01]; a National Science Foundation INSPIRE award [grant number 1649014]; and the National Science Foundation / ERASynBio [grant numbers 1548123, 335672].

1.4.4 Questions and comments

Please contact the [Karr Lab](#) with any questions or comments.

1.5 References